
DATA MINING V

DATA MINING, TEXT MINING
AND THEIR BUSINESS APPLICATIONS

WIT*PRESS*

WIT Press publishes leading books in Science and Technology.

Visit our website for the current list of titles.

www.witpress.com

WIT*eLibrary*

Home of the Transactions of the Wessex Institute, the WIT electronic-library provides the international scientific community with immediate and

permanent access to individual papers presented at WIT conferences.

Visit the WIT eLibrary at www.witpress.com

FIFTH INTERNATIONAL CONFERENCE ON DATA MINING
DATA MINING V

CONFERENCE CHAIRMEN

A. Zanasi

TEMIS Text Mining Solutions, Italy

N.F.F. Ebecken

Federal University of Rio de Janeiro, Brazil

C.A. Brebbia

Wessex Institute of Technology, UK

INTERNATIONAL SCIENTIFIC ADVISORY COMMITTEE

- S. Ananyan *Megaputer, Inc., USA*
M. Berry *University of Tennessee, USA*
S. Bolasco *University of Rome, Italy*
R. Maspons Bosch *IALE Technologia, Spain*
O. Ciftcioglu *Delft Univ. of Technology, Netherlands*
M. Costantino *U.K.*
P. Coupet *TEMIS SA, France*
B. Drewes *SAS Institute, Germany*
P. Giudici *University of Pavia, Italy*
T. Khabaza *SPSS, U.K.*
G. Lachtermacher *Faculdades IBMEC, Brazil*
D. Laney *META Group, USA*
A. Linden *Gartner Group, Germany*
D. Malerba *Universita degli Studi, Italy*
G. Marchisio *Insightful, U.S.A.*
H. Messatfa *IBM Consulting Services, France*
N.M. Milic-Frayling *Microsoft Research Ltd, U.K.*
E. Orozco *IDICT, Cuba*
P. J.-S. Pan *Nat Kaohsiung Uni. of Ap. Sci., Taiwan*
M.F.R. Rodrigues *Poly Institute of Porto, Portugal*
O. Ryabov *Kharkov Academy of Culture, Ukraine*
D. Sacca *DEIS, Italy*
S. Sirmakessis *University of Patras, Greece*
D. Sitnikov *Kharkov Academy of Culture, Ukraine*
R. Turra *CINECA, Italy*
D.E.N. Van den Poel *Ghent University, Belgium*
R. Weber *Universidad de Chile, Chile*
N. Zhong *Maebashi Inst.of Technology, Japan*
H.G. Zimmermann *Siemens AG, Germany*

Organised by

Wessex Institute of Technology, UK

DATA MINING V

DATA MINING, TEXT MINING
AND THEIR BUSINESS APPLICATIONS

Editors

A. Zanasi

TEMIS Text Mining Solutions, Italy

N.F.F. Ebecken

Federal University of Rio de Janeiro, Brazil

C.A. Brebbia

Wessex Institute of Technology, UK

WITPRESS Southampton, Boston



A. Zanasi

TEMIS Text Mining Solutions, Italy

N.F.F. Ebecken

Federal University of Rio de Janeiro, Brazil

C.A. Brebbia

Wessex Institute of Technology, UK

Published by

WIT Press

Ashurst Lodge, Ashurst, Southampton, SO40 7AA, UK

Tel: 44 (0) 238 029 3223; Fax: 44 (0) 238 029 2853

E-Mail: witpress@witpress.com

<http://www.witpress.com>

For USA, Canada and Mexico

WIT Press

25 Bridge Street, Billerica, MA 01821, USA

Tel: 978 667 5841; Fax: 978 667 7582

E-Mail: infousa@witpress.com

<http://www.witpress.com>

British Library Cataloguing-in-Publication Data

A Catalogue record for this book is available
from the British Library

ISBN: 1-85312-722-9

ISSN: 1470-6326

*The texts of the papers in this volume were set
individually by the authors or under their supervision.
Only minor corrections to the text may have been carried
out by the publisher.*

No responsibility is assumed by the Publisher, the Editors and Authors for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions or ideas contained in the material herein.

© WIT Press 2004.

Printed in Great Britain by Athenaeum Press Ltd. Gateshead.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the Publisher.

Preface

The Data Mining Conference in 2004 was again held in Europe, in Malaga, Spain, after previously taken place in Rio de Janeiro, Brazil, which hosted the 4th Conference.

This book contains the papers presented at the 5th Data Mining Conference. The Conference is becoming a traditional meeting for the practitioners of data mining. As data mining is evolving, so is the conference, although it remains faithful to its original formula of providing a meeting point for practitioners of data mining from academia and industry, as can be seen in the contributions published in this book.

Nearly 25 countries from every continent are represented in this book. The ISAC-International Scientific Advisory Committee includes high level scientists and professionals representing academia (from North and South America, Europe and Asia), industries (IBM, Microsoft, Siemens, SAS, SPSS, Insightful, Megaputer, Temis) and analysts (Gartner Group and META Group).

Data Mining is evolving, but in what sense and in what direction? The definition of data mining as “the process of extracting previously unknown, valid and actionable information from large databases and then using it to make crucial business decisions”¹ is still valid in the sense that many other sectors, such as artificial intelligence, applied mathematics, control theory and information technology still maintain an important role in providing techniques and subjects of research interesting to data mining.

It is also important that the concept of “large databases” be understood in a wider meaning by taking into account the, so called, “internet sources.” Until recently mining has been, in fact, relegated to well-structured internal information. Mining the public web has been an exercise in frustration, but it is now possible to apply the same mining tools and techniques to the explosion of data found on billions of web pages, newsgroups, emails, chatlines, forums, blogs and message boards².

¹ Cabena, Hadjinian, Stadler, Verhees, Zanasi-*Discovering Data Mining*-1998-Prentice Hall

² A.Zanasi-Email, chatlines, newsgroups: a continuous opinion surveys source thanks to text mining application – 2003 –in *Excellence 2003 in Int'l Research* - Ed.Esomar

In this sense, thanks to the explosion of the Internet phenomenon and of new sources, data mining applied to text, ie. text mining, has acquired much more importance.

Researchers coming from other disciplines of science and research (eg. information retrieval, semantic web, linguistics, knowledge management...) are starting to show interest in *unstructured data mining* (ie text mining) results for their studies and applications.

This is the reason why it was decided to add a clear reference to text mining in the title of this conference, which now appears as “Data Mining, Text Mining and its Business Applications”, to be retained also for the next conference (Data Mining 2005, to be held on 25 May 2005 in Skiathos, Greece).

The reference to “Business Applications” is there to stress that contributions regarding solutions to the problems that the real world pose to us are particularly welcome.

In recent years, businesses operating in diverse areas have shown an interest in data mining. Three of them appear to be the most promising, i.e.

National Security. Reading magazines, listening to television all over the world, no matter where one lives, there is no doubt that the most pressing problems are those linked to terrorism and to national security.

It is not a coincidence that more books are being published on data, text mining and business intelligence techniques applied to intelligence or national or homeland security.

In the USA fierce discussions about the opportunity of fighting terrorism using data and text mining capabilities, organized under DARPA (eg Information Awareness Office and Terrorist Information Awareness programs), are still firing the political debate, fearing the risk that this means to data privacy and to democracy itself. The CIA launched In-Q-Tel, a non profit corporation designed to fuel private research, development and application in, amongst others, data and text mining technologies³. It is to be expected that in the next few years this business area will grow quickly.

Competitive Intelligence (CI). As one of the developing software-based new analysis techniques, data and text mining provide a powerful way to uncover changing business issues and discover emerging trends that can drive business. CI capabilities are also changing as professionals pick up the tools that can help them from being reporters of facts, to predictors of market changes.

Combining the precision of analysis with the serendipity of discovery, text mining helps detect the complex signals of emerging trends over billions of pieces of information. Together with visualization tools, text mining provides a unique window to the world⁴.

³ John Gannon-*The Strategic Use of Open-Source Information* –in Intelligence and the National Security Strategist: Enduring Issues and Challenges-Sherman Kent Center for Intelligence Studies National War College National Defense University (2004)

⁴ Dennis Cahill – Visualizing Emerging Intelligence through Text Mining – in Competitive Intelligence Magazine, May-June 2004

Customer Relationship Management (CRM). Many organizations already apply data mining to internal customer information. Using large Customer Relationship Management (CRM) tools and other internal databases, companies identify in a detailed way how customers purchase (or don't purchase) products. In the high-volume consumer transactions of ecommerce business, analyzing information through data mining can yield intelligence that can affect revenue or profit. The opportunity of analyzing emails and chatlines, thanks to text mining, allow for the detection of customers' opinions and trends to boost sales. It is not a coincidence that a section of this book is dedicated to CRM.

For a discipline that is central to research and business, and which is of interests not only among software developers, but also among large companies and government departments and agencies, it is important to define its market size. IDC estimate, for only the data mining market size in 2005⁵, \$ 1.5 billion, and OVUM estimate for data and text mining and the related technologies, grouped under the term "knowledge management"⁶, \$ 11.98 billion. These are large sums, and seeing the trend of growing government expenses in the national security market, where data and text mining are of fundamental importance, this sum will continue to grow.

The success of the Conference and this resulting book would not have been possible without the generous help of the members of the I.S.A.C. Their efforts ensured the quality of the papers published in this volume.

Alessandro ZANASI

Malaga, September 2004

⁵ IDC-Information Access Tools : Market Forecast and Analysis

⁶ SHIC Journal: <http://www.siic.org.hk/journal/935.582.107.PDF>

Contents

SECTION 1: TEXT MINING

Personalization in the semantic web era: a glance ahead <i>P. Markellou, M. Rigou, S. Sirmakessis & A. Tsakalidis</i>	3
A simple mixture model for unsupervised text categorisation <i>F. Clérot, F. Fessant, O. Collin, O. Cappé & E. Moulines</i>	13
Data warehouse screening and personalizing agent <i>A. H. Mohamed, N. L. M. Noor & N. M. Noh</i>	23
Extracting people's names from RSS feeds using WordNet <i>K. Durant & M. Smith</i>	33
Stock Broker P – sentiment extraction for the stock market <i>R. Khare, N. Pathak, S. K. Gupta & S. Sohi</i>	43
An XML based semantic protein map <i>A. S. Sidhu, T. S. Dillon & H. Setiawan</i>	51
Automated text mining comparison of Japanese and USA multi-robot research <i>R. J. Watts, A. Porter & B. Minsk</i>	61

SECTION 2: WEB MINING

The influence of caching on web usage mining <i>J. Huysmans, B. Baesens & J. Vanthienen</i>	77
Clickstreams, the basis to establish user navigation patterns on web sites <i>R. Alves, O. Belo, F. Cavalcanti & P. Ferreira</i>	87

SECTION 3: TECHNIQUES

Section 3a: Clustering

Robust clustering methods for incomplete and erroneous data <i>T. Kärkkäinen & S. Áyrämö</i>	101
Clustering as an add-on for firewalls <i>C. Caruso & D. Malerba</i>	113
Exploration of the ecological status of Mediterranean rivers: clustering, visualizing and reconstructing streams data using Generative Topographic Mapping <i>D. Vicente, A. Vellido, E. Martí, J. Comas & I. Rodriguez-Roda</i>	121

Section 3b: Categorization

A fuzzy decision tree approach to start a genetic algorithm for data classification <i>R. P. Espindola & N. F. F. Ebecken</i>	133
Fuzziness as a recognition problem: using decision tree learning algorithms for inducing fuzzy membership functions <i>O. Nykänen</i>	143

Section 3c: Link analysis

Mining association rules with negative terms using candidate pruning <i>T. Shintani & D. Hayashi</i>	157
Association rules model of e-banking services <i>V. Aggelis</i>	167

Section 3d: Data preparation

An algebraic approach to defining rough set approximations and generating logic rules <i>D. Sitnikov & O. Ryabov</i>	179
Scalability issue in mining large data sets <i>A. Mc Manus & M.-T. Kechadi</i>	189
A Bayesian approach for supervised discretization <i>M. Boullé</i>	199

A statistical and signal processing based system for data quality management <i>A. C. H. Dantas, J. M. de Seixas, F. B. Diniz & T. N. Ferreira</i>	209
---	-----

SECTION 4: APPLICATIONS IN BUSINESS, INDUSTRY AND GOVERNMENT

An e-Knowledge application for local government and small and medium businesses <i>M. Castellano, N. Pastore, F. Arcieri, V. Summo & G. Bellone de Grecis</i>	221
--	-----

Daily sugar price forecasting using the Mixture of Local Expert Models <i>B. de Melo, C. L. Nascimento Júnior & A. Z. Milioni</i>	231
--	-----

A visual tool for mining macroeconomics data <i>D. Giordano & F. Maiorana</i>	241
--	-----

Data mining in publishing: a nice feature or a necessity? <i>T. M. Fernández-Steeger, F. Zander, S. Callsen, S. Steinberg & N. Brauns</i>	253
--	-----

A practical implementation of a real-time intrusion prevention system for commercial enterprise databases <i>U. T. Mattsson</i>	263
--	-----

The applications of genetic algorithms in stock market data mining optimisation <i>L. Lin, L. Cao, J. Wang & C. Zhang</i>	273
--	-----

Genetic algorithm applied on the performance appraisal system of mutual fund managers in Taiwan <i>J.-F. Chang & B.-Y. Liao</i>	281
--	-----

Collusion in the US crop insurance program: applied data mining <i>B. B. Little, W. L. Johnston, Jr., A. C. Lovell, R. M. Rejesus & S. A. Steed</i>	291
--	-----

SECTION 5: CUSTOMER RELATIONSHIP MANAGEMENT

Analysis of service quality indicators in telecommunications using control statistical methods and fuzzy logic <i>C. A. A. Lemos, N. F. F. Ebecken & A. G. Evsukoff</i>	305
--	-----

Mining call center dialog data <i>A. Gilman, B. Narayanan & S. Paul</i>	317
A naive KDD approach in a Key Account Management Framework: a case study <i>C. A. F. Gama, A. Evsukoff & J. P. Motta</i>	327
Customer valuation through data mining <i>C. Fernando Nogueira & N. F. F. Ebecken</i>	333
A data mining approach to analysis and prediction of movie ratings <i>M. Saraee, S. White & J. Eccleston</i>	343
 SECTION 6: APPLICATIONS IN SCIENCE AND ENGINEERING	
New search strategies and a new derived inequality for efficient k-medoids-based algorithms <i>C.-S. Chiang, S.-C. Chu, J. F. Chang & J.-S. Pan</i>	355
A data mining approach to support the development of new fuels and technology <i>G. S. Terra, C. L. Curotto & N. F. F. Ebecken</i>	365
Data mining highly multiple time series of astronomical observations <i>F. Huang</i>	375
A Model-View-Controller architecture for Knowledge Discovery <i>M. Castellano, N. Pastore, F. Arcieri, V. Summo & G. Bellone de Grecis</i>	383
MZ-Platform: a component switching and executing environment <i>N. Matsuki</i>	393
Data mining approach to study Quality of Voice over IP applications <i>I. Miloucheva, D. Hetzer & A. Nassri</i>	403
Data mining and population genetics of birth defects: preliminary investigation <i>B. Little</i>	415
A data mining approach to landslide prediction <i>F. T. Souza & N. F. F. Ebecken</i>	423
 Author Index	 433